

IN THE UNITED STATES PATENT AND TRADEMARK OFFICE

APPLICATION FOR LETTERS PATENT

**Locating Potentially Identical Objects Across Multiple  
Computers Based On Stochastic Partitioning Of  
Workload**

Inventor(s):

**John R. Douceur**  
**Marvin M. Theimer**  
**Atul Adya**  
**William J. Bolosky**

ATTORNEY'S DOCKET NO. MS1-769US

## **TECHNICAL FIELD**

This invention relates to computer networks and file systems, and more particularly to locating potentially identical files across multiple computers based on stochastic partitioning of workload.

## **BACKGROUND OF THE INVENTION**

File systems manage files and other data objects stored on computer systems. File systems were originally built into the computer operating system to facilitate access to files stored locally on resident storage media. As personal computers became networked, some file storage capabilities were offloaded from individual user machines to special storage servers that stored large numbers of files on behalf of the user machines. When a file was needed, the user machine simply requested the file from the server. In this server-based architecture, the file system is extended to facilitate management of and access to files stored remotely at the storage server over a network.

One problem that arises in distributed file systems concerns storage of identical files on the servers. While some file duplication normally occurs on an individual user's personal computer, duplication unfortunately tends to be quite prevalent on networks where servers centrally store the contents of multiple personal computers. For example, with a remote boot facility on a computer network, each user boots from that user's private directory on a file server. Each private directory thus ordinarily includes a number of files that are identical to files on other users' directories. Storing the private directories on traditional file systems consumes a great amount of disk and server file buffer cache space. From a storage management perspective, it is desirable to reduce file duplication to

1 reduce the amount of wasted storage space used to store redundant files.  
2 However, any such efforts need to be reconciled with the file system that tracks  
3 the multiple duplicated files on behalf of the associated users.

4 To address the problems associated with storing multiple identical files on a  
5 computer, Microsoft developed a single instance store (SIS) system that is  
6 packaged as part of the Windows 2000 operating system. The SIS system reduces  
7 file duplication by automatically identifying common identical files of a file  
8 system, and then merging the files into a single instance of the data. One or more  
9 logically separate links are then attached to the single instance to represent the  
10 original files to the user machines. In this way, the storage impact of duplicate  
11 files on a computer system is greatly reduced.

12 Today, file storage is migrating toward a model in which files are stored on  
13 various networked computers, rather than on a central storage server. However,  
14 the problem of duplicate identical files remains, except that the duplicate files are  
15 spread out over the various networked computers. Given the large number of  
16 computers that can currently be networked together (easily into the thousands or  
17 hundreds of thousands), and the large number of files that can exist spread out  
18 over this large number of computers (easily into the millions or billions), detecting  
19 duplicate files in such an environment can be very difficult. Limitations on the  
20 bandwidth available to transfer information among the computers, as well as  
21 limitations on the computational capacity of the computers themselves, makes  
22 such detections very difficult.

23 The invention addresses these problems, allowing locating of potentially  
24 identical objects, such as files, across multiple computers.  
25

## SUMMARY OF THE INVENTION

Locating potentially identical objects across multiple computers based on stochastic partitioning of workload is described herein.

In accordance with one aspect, identical objects (e.g., files) are located across multiple computers by selecting, for each of a plurality of objects stored on a plurality of computers in a network, a portion of object information corresponding to the object. The object information can be generated in a variety of manners (e.g., based on hashing the object, based on characteristics of the object, and so forth). Any of a variety of portions of the object information can be used (e.g., the least significant bits of the object information). A stochastic partitioning process is then used to identify which of the plurality of computers to communicate the object information to for identification of potentially identical objects on the plurality of computers.

According to another aspect, the stochastic partitioning process comprises a fully distributed stochastic partitioning process including in which, for each of a plurality of computers, the selected portion of the object information is compared to a portion of a computer identifier associated with the computer. An identification is then made as to which of the computer identifiers have portions matching the selected portion of the object information, and the object information is communicated to each of the computers associated with a computer identifier having a portion matching the selected portion of the object information.

According to another aspect, the stochastic partitioning process comprises a group-based system using directory services process in which an object information portion to computer mapping is accessed on a remote computer.

1 Based on the selected portion of the object information and the remotely accessed  
2 mapping, one or more computers are identified to receive the object information.

3 According to another aspect, the stochastic partitioning process comprises a  
4 stochastic partitioning process comprises a multi-level stochastic partitioning  
5 process in which selected ones of the plurality of computers in the network are  
6 grouped into a plurality of groups based at least in part on the number of the  
7 plurality of computers in the network that the computer using the stochastic  
8 partitioning process is aware of. Which of the selected ones of the plurality of  
9 computers to communicate the object information to is then identified, wherein the  
10 identifying is based at least in part on comparing the selected portion of the object  
11 information to a portion of a computer identifier of one or more of the selected  
12 ones of the plurality of computers.

### 13 14 **BRIEF DESCRIPTION OF THE DRAWINGS**

15 The present invention is illustrated by way of example and not limitation in  
16 the figures of the accompanying drawings. The same numbers are used  
17 throughout the figures to reference like components and/or features.

18 Fig. 1 illustrates an exemplary network environment that supports a  
19 serverless distributed file system.

20 Fig. 2 illustrates logical components of an exemplary computing device that  
21 is representative of any one of the devices of Fig. 1 that participate in the  
22 distributed file system.

23 Fig. 3 is a flowchart illustrating an exemplary process followed to inform  
24 database servers of the file information in accordance with certain embodiments of  
25 the invention.

1 Fig. 4 is a flowchart illustrating an exemplary process followed by a  
2 database server in accordance with certain embodiments of the invention.

3 Fig. 5 illustrates an exemplary centralized database implementation in  
4 accordance with certain embodiments of the invention.

5 Fig. 6 illustrates an exemplary network environment in which responsibility  
6 for managing the file information database is shared among multiple database  
7 servers in each group.

8 Fig. 7 illustrates a localized pair-wise checking implementation in  
9 additional detail.

10 Fig. 8 illustrates the special situation when the number of client computers  
11 in a group is equal to all of the computers in the network in additional detail.

12 Fig. 9 is a flowchart illustrating an exemplary process followed by each  
13 computer for the fully distributed stochastically partitioned database  
14 implementation in accordance with certain embodiments of the invention.

15 Fig. 10 illustrates an exemplary network in which a fully distributed  
16 stochastically partitioned database implementation is employed.

17 Fig. 11 is a flowchart illustrating an exemplary process followed by each  
18 computer for the group-based system using directory services implementation in  
19 accordance with certain embodiments of the invention.

20 Fig. 12 illustrates an exemplary network in which the group-based system  
21 using directory services implementation is employed.

22 Fig. 13 is a flowchart illustrating an exemplary process followed by each  
23 computer for a multi-level stochastically partitioned database implementation in  
24 accordance with certain embodiments of the invention.  
25

1 Fig. 14 illustrates an exemplary network in which a multi-level  
2 stochastically partitioned database implementation is employed.

3 Fig. 15 illustrates a more general exemplary computer environment which  
4 can be used in various embodiments of the invention.

## 5 6 **DETAILED DESCRIPTION**

7 The discussions herein assume a basic understanding of cryptography by  
8 the reader. For a basic introduction of cryptography, the reader is directed to a  
9 text written by Bruce Schneier and entitled "Applied Cryptography: Protocols,  
10 Algorithms, and Source Code in C," published by John Wiley & Sons with  
11 copyright 1994 (or second edition with copyright 1996).

### 12 13 **Operating Environment**

14 The following discussion is directed primarily to locating identical files  
15 across multiple computers in a distributed file system. The file system is  
16 described in the context of a symbiotic, serverless, distributed file system that runs  
17 on multiple networked computers and stores files across the computers rather than  
18 on a central server or cluster of servers. The symbiotic nature implies that the  
19 machines cooperate but do not completely trust one another. The file system does  
20 not manage the storage disk directly, but rather relies on existing file systems on  
21 local machines, such as those file systems integrated into operating systems (e.g.,  
22 the Windows NT® file system).

23 While the file system is described in the context of storing "files", it should  
24 be noted that other types of storable data can be stored in the file system. The  
25 term "file" is used for discussion purposes and is intended to include data objects

1 or essentially any other storage subject matter that may not be commonly  
2 characterized as a "file".

3 Additionally, the systems and methods described herein are also applicable  
4 to data in other types of systems other than file systems, such as database systems  
5 or object systems. The methods and systems described herein operate on objects  
6 containing bytes (these objects being predominately described herein as files), and  
7 can be used to identify potentially duplicate objects including any object data or  
8 meta data. Furthermore, the methods and systems described herein may also  
9 operate on object-defined methods rather than acting upon the objects at the byte  
10 level, including methods defined on objects for obtaining the bytes (e.g., file  
11 information) described herein.

12 Fig. 1 illustrates an exemplary network environment 100 that supports a  
13 serverless distributed file system. Four client computing devices 102, 104, 106,  
14 and 108 are coupled together via a data communications network 110. Although  
15 four computing devices are illustrated, different numbers (either greater or fewer  
16 than four) may be included in network environment 100.

17 Network 110 represents any of a wide variety of data communications  
18 networks. Network 110 may include public portions (e.g., the Internet) as well as  
19 private portions (e.g., an internal corporate Local Area Network (LAN)), as well  
20 as combinations of public and private portions. Network 110 may be implemented  
21 using any one or more of a wide variety of conventional communications media  
22 including both wired and wireless media. Any of a wide variety of  
23 communications protocols can be used to communicate data via network 110,  
24 including both public and proprietary protocols. Examples of such protocols  
25 include TCP/IP, IPX/SPX, NetBEUI, etc.



1 Computing devices 102-108 represent any of a wide range of computing  
2 devices, and each device may be the same or different. By way of example,  
3 devices 102-108 may be desktop computers, laptop computers, handheld or pocket  
4 computers, personal digital assistants (PDAs), cellular phones, Internet appliances,  
5 consumer electronics devices, gaming consoles, and so forth.

6 Two or more of devices 102-108 operate to implement a serverless  
7 distributed file system (although some of them may not be operational (e.g., failed  
8 or powered-down) at any given time). The actual devices included in the  
9 serverless distributed file system can change over time, allowing new devices to  
10 be added to the system and other devices to be removed from the system. Each  
11 device 102-108 that is part of the distributed file system has different portions of  
12 its mass storage device(s) (e.g., hard disk drive) allocated for use as either local  
13 storage or distributed storage. The local storage is used for data that the user  
14 desires to store on his or her local machine and not in the distributed file system  
15 structure. The distributed storage portion is used for data that the user of the  
16 device (or another device) desires to store within the distributed file system  
17 structure.

18 In the illustrated example of Fig. 1, certain devices connected to network  
19 110 have one or more mass storage devices that include both a portion used by the  
20 local machine and a portion used by the distributed file system. The amount  
21 allocated to distributed or local storage varies among the devices and can vary  
22 over time. For example, device 102 has a larger percentage allocated for a  
23 distributed system portion 120 in comparison to the local portion 122; device 104  
24 includes a distributed system portion 124 that is approximately the same size as  
25 the local portion 126; and device 106 has a smaller percentage allocated for a

1 distributed system portion 128 in comparison to the local portion 130. The storage  
2 separation into multiple portions may occur on a per storage device basis (e.g., one  
3 hard drive is designated for use in the distributed system while another is  
4 designated solely for local use), and/or within a single storage device (e.g., part of  
5 one hard drive may be designated for use in the distributed system while another  
6 part is designated for local use). Other devices connected to network 110, such as  
7 computing device 108, may not be part of the distributed file system and thus such  
8 devices do not have any of their mass storage device(s) allocated for use by the  
9 distributed system. Hence, device 108 has only a local portion 132.

10 A distributed file system 150 operates to store one or more copies of files  
11 on different computing devices 102-106. When a new file is created by the user of  
12 a computer, he or she has the option of storing the file on the local portion of his  
13 or her computing device, or alternatively in the distributed file system. If the file  
14 is stored in the distributed file system 150, the file will be stored in the distributed  
15 system portion of the mass storage device(s) of one or more of devices 102-106.  
16 The user creating the file typically has no ability to control which device 102-106  
17 the file is stored on, nor any knowledge of which device 102-106 the file is stored  
18 on. Additionally, replicated copies of the file will typically be saved, allowing the  
19 user to subsequently retrieve the file even if one of the computing devices 102-106  
20 on which the file is saved is unavailable (e.g., is powered-down, is malfunctioning,  
21 etc.).

22 The distributed file system 150 is implemented by one or more components  
23 on each of the devices 102-106, thereby obviating the need for any centralized  
24 server to coordinate the file system. These components operate to determine  
25 where particular files are stored, how many copies of the files are created for

1 storage on different devices, and so forth. Exactly which device will store which  
2 files depends on numerous factors, including the number of devices in the  
3 distributed file system, the storage space allocated to the file system from each of  
4 the devices, how many copies of the file are to be saved, the number of files  
5 already stored on the devices, and so on. Thus, the distributed file system allows  
6 the user to create and access files (as well as folders or directories) without any  
7 knowledge of exactly which other computing device(s) the file is being stored on.

8 The files stored by the file system are distributed among the various devices  
9 102-106 and stored in encrypted form. When a new file is created, the device on  
10 which the file is being created encrypts the file prior to communicating the file to  
11 other device(s) for storage. The directory entry (e.g., the file name) for a new file  
12 is also communicated to the other device(s) for storage. Additionally, if a new  
13 folder or directory is created, the directory entry (e.g., folder name or directory  
14 name) is also communicated to the other device(s) for storage. As used herein, a  
15 directory entry refers to any entry that can be added to a file system directory,  
16 including both file names and directory (or folder) names.

17 The distributed file system 150 is designed to prevent unauthorized users  
18 from reading data stored on one of the devices 102-106. Thus, a file created by  
19 device 102 and stored on device 104 is not readable by the user of device 104  
20 (unless he or she is authorized to do so). In order to implement such security, the  
21 contents of files as well as all directory entries are encrypted, and only authorized  
22 users are given the decryption key. Thus, although device 104 may store a file  
23 created by device 102, if the user of device 104 is not an authorized user of the  
24 file, the user of device 104 cannot decrypt (and thus cannot read) either the  
25 contents of the file or its directory entry (e.g., filename).

Fig. 2 illustrates logical components of an exemplary computing device 200 (also referred to herein as a computer or machine) that is representative of any one of the devices 102-106 of Fig. 1 that participate in the distributed file system 150. Computing device 200 includes a mass storage device 208, a distributed file system interface 210, and various additional modules providing client and/or server functionality. Computing device 200 also typically includes additional components (e.g., a processor), however these additional components have not been shown in Fig. 2 so as not to clutter the drawings. A more general description of a computer architecture with various hardware and software components is described below with reference to Fig. 15.

Mass storage device 208 can be any of a wide variety of conventional nonvolatile storage devices, such as a magnetic disk, optical disk, Flash memory, and so forth. Mass storage device 208 is separated into a distributed system portion and a local portion; this separation may change over time.

Computing device 200 is intended to be used in a serverless distributed file system, and as such includes modules oriented towards both server functionality and client functionality. The server functionality comes into play when device 200 is responding to a request involving a file or directory entry stored (or to be stored) in storage device 208, as well as when identifying potentially identical or duplicate files. The client functionality, on the other hand, comes into play when issuing requests by device 200 for files stored (or to be stored) in the distributed file system, as well as generating and forwarding file information for file duplication identification as necessary. The client and server functionality operate independent of one another. Thus, situations can arise where the serverless distributed file system 150 causes files being stored by modules operating in a

1 client capacity to be stored in mass storage device 208 by other modules operating  
2 in a server capacity.

3 Computing device 200 includes a file information generation module 220, a  
4 forwarding location determination module 222, and a file information comparison  
5 module 242. It should be noted, however, that not all components are necessarily  
6 needed on each computing device 200. For example, a computing device  
7 operating as a client-only machine might not include file information comparison  
8 module 242, or a computing device operating as a server-only machine might not  
9 include file information generation module 220.

10 File information generation module 220 generates file information for one  
11 or more of encrypted files 240 in storage device 208. Forwarding location  
12 determination module 222 determines the location (e.g., one or more other  
13 computing devices 200) where the file information generated by module 222 is to  
14 be communicated. These components and their operation are described in more  
15 detail below.

16 Although illustrated on a single computing device in Fig. 2, file information  
17 generation module 220 and forwarding location determination module 222 may  
18 also be implemented across multiple computing devices. For example, in the  
19 distributed file system environment illustrated in Fig. 1, a computing device may  
20 create or update a file for storage in distributed file system 150, and then  
21 communicate the file to another device(s) in distributed file system 150 acting as a  
22 directory server. The directory server then stores the file on an appropriate  
23 computing device (based on the rules followed by distributed file system 150) and  
24 maintains a record of where the file is stored. In this example, the computing  
25 device creating or updating the file generates the file information (via its file

1 information generation module 220), while the computing device acting as the  
2 directory server (and thus which knows what other computing device the file is  
3 stored at) determines the location where the generated file information is to be  
4 communicated (via its forwarding location determination module 222).

5 File information generation module 220 generates file information for one  
6 or more encrypted files 240. The file information for each file is a semi-unique  
7 value based on the data in the file itself (the data may be program instructions,  
8 program data, etc.) and/or other characteristics of the file. The value is a semi-  
9 unique value because it is based on the data in the file but is not completely  
10 representative of the file. For example, the file information may be a hash value  
11 that is based on the data in the file, but it is possible for two different files having  
12 different data to have the same hash value. Different characteristics of the file can  
13 also be incorporated into the file information, such as the file size, the file type, the  
14 file name, and so forth. The file information can be generated in any of a wide  
15 variety of manners, so long as each of the computing devices generates its file  
16 information in the same manner. Two files that have different file information are  
17 not duplicate files. Two files that have the same file information, however, may or  
18 may not be duplicate files.

19 In one implementation, the file information is a hash value generated based  
20 on the file. The hash value may be generated using a one-way hashing function  
21 (e.g., SHA, MD5, etc.), or any of a variety of other public or proprietary hashing  
22 functions. The hash value may be based on the entire file, or alternatively only a  
23 portion of the file (e.g., the beginning of the file, the end of the file, the middle of  
24 the file, and so forth). In another implementation, the file information is referred  
25 to as a file signature, which is a combination of a hash value based on the file (the

1 hash value represents 64 bits of the file signature) and the file size (which  
2 represents another 64 bits of the file signature). The file size is used because two  
3 files with differing file sizes cannot be identical.

4 In certain embodiments, the file information is based on a hash value  
5 corresponding to the file. This hash value is based on both block-by-block  
6 encryption and convergent encryption (as described below), and is generated by  
7 the file system for secure storage of files in the distributed computing  
8 environment. Thus, in these embodiments, file information generation module  
9 220 need only access the hash value already generated for a file for secure storage  
10 in order to generate the file information.

11 Generally, according to convergent encryption, a file  $F$  (or any other type of  
12 encryptable object) is initially hashed using a one-way hashing function  $h$  (e.g.,  
13 SHA, MD5, etc.) to produce a hash value  $h(F)$ . The file  $F$  is then encrypted using  
14 a symmetric cipher (e.g., RC4, RC2, etc.) with the hash value as the key, or  
15  $E_{h(F)}(F)$ . Next, read access control entries are created for each authorized user who  
16 is granted read access to the encrypted file. The access control entries are formed  
17 by encrypting the file's hash value  $h(F)$  with any number of keys  $K_1, K_2, \dots, K_m$ ,  
18 to yield  $E_{K_1}(h(F)), E_{K_2}(h(F)), \dots, E_{K_m}(h(F))$ . The keys  $K$  are randomly generated  
19 and uniquely assigned to individual users. In one implementation, each key  $K$  is  
20 the user's public key of a public/private key pair. In the illustrated example, write  
21 access control is governed by the directory server that stores the directory entry for  
22 the file and it is thus not addressed by the file format (so references to "access"  
23 within this document refer to read access unless specifically identified as another  
24 type of access). Alternatively, write access control could be implemented via  
25

1 access control entries in a manner analogous to the read access control discussed  
2 herein.

3 With convergent encryption, one encrypted version of the file is stored and  
4 replicated among the serverless distributed file system 150. Along with the  
5 encrypted version of the file is stored one or more access control entries depending  
6 upon the number of authorized users who have access. Thus, a file in the  
7 distributed file system 150 has the following structure:

$$8 \quad [E_{h(F)}(F), \langle E_{K1}(h(F)) \rangle, \langle E_{K2}(h(F)) \rangle, \dots, \langle E_{Km}(h(F)) \rangle]$$

11 One advantage of convergent encryption is that the encrypted file can be  
12 evaluated by the file system to determine whether it is identical to another file  
13 without resorting to any decryption (and hence, without knowledge of any  
14 encryption keys). Unwanted duplicative files can be removed by adding the  
15 authorized user(s) access control entries to the remaining file. Another advantage  
16 is that the access control entries are very small in size, on the order of bytes as  
17 compared to possibly gigabytes for the encrypted file. As a result, the amount of  
18 overhead information that is stored in each file is reduced. This enables the  
19 property that the total space used to store the file is proportional to the space that is  
20 required to store a single encrypted file, plus a constant amount of storage for each  
21 additional authorized reader of the file.

22 For more information on convergent encryption, the reader is directed to  
23 co-pending U.S. Patent Application Serial No. 09/565,821, entitled "Encryption  
24 Systems and Methods for Identifying and Coalescing Identical Objects Encrypted  
25 with Different Keys", which was filed May 5, 2000, in the names of Douceur et



1 al., and is commonly assigned to Microsoft Corporation. This application is  
2 hereby incorporated by reference.

3 For small files, the entire file is hashed and encrypted using convergent  
4 encryption, and the resulting hash value is used as the encryption key. The  
5 encrypted file can be verified without knowledge of the key or any need to decrypt  
6 the file first. For large files, the file contents are broken into smaller blocks and  
7 then convergent encryption is applied separately to each block. For example, the  
8 file  $F$  may be segmented into "n" pages  $F^0$ - $F^{n-1}$ , where each page is a fixed size  
9 (e.g., a 4Kbyte size). Convergent encryption is then applied to the file at the block  
10 level. That is, each block  $F^i$  is separately hashed using a one-way hash function  
11 (e.g., SHA, MD5, etc.) to produce a hash value  $h(F^i)$ . Each block  $F^i$  is then  
12 encrypted using a symmetric cipher (e.g., RC4, RC2, etc.) with the hash value  
13  $h(F^i)$  as the key, or  $E_{h(F^i)}(F^i)$ , resulting in an array of encrypted blocks which form  
14 the contents of the file. For more information on block-by-block encryption, the  
15 reader is directed to co-pending U.S. Patent Application Serial No. \_\_\_\_\_ entitled  
16 "On-Disk File Format for Serverless Distributed File System", Attorney Docket  
17 No. MS1-733US, to inventors William J. Bolosky, Gerald Cermak, Atul Adya, and  
18 John R. Douceur. This application is hereby incorporated by reference.

19 File information generation module 220 can generate the file information at  
20 any of a wide variety of times. In one implementation, module 220 is designed to  
21 operate as a background process. When files are created or modified, the file  
22 names are added to a queue to be acted on by module 220. When computing  
23 device 200 is not busy (e.g., the processor has free cycles, or has been idle for a  
24 period of time), module 220 operates to generate file information for one of the  
25 files in the queue. Alternatively, module 220 may be designed to run at times of

1 anticipated low usage (e.g., at night or early morning), or module 220 may  
2 generate the file information for a file whenever that file is created or modified.

3 Module 220 may generate file information for each encrypted file 240, or  
4 alternatively only for selected files 240. In one implementation, module 220  
5 generates file information only for files greater than a threshold size (e.g., files that  
6 are at least 16k bytes). This threshold size is implemented to account for the  
7 situation where the overhead necessary to identify and coalesce duplicate files that  
8 are very small is deemed to be too great in light of the small amount of storage  
9 space (due to the small file size) that could be recovered.

10 The file duplication identification described herein is described primarily  
11 with reference to files 240 stored in the distributed system portion(s) of storage  
12 device 208. Alternatively, the file duplication identification could also be applied  
13 to files stored in the local portion(s) of storage device 208.

14 Forwarding location determination module 222 receives the file  
15 information from file information generation module 220 and forwards the file  
16 information to one or more other computing devices 200. Which other computing  
17 devices the file information is forwarded to can vary, and is discussed in more  
18 detail below with respect to the various implementations.

19 Additionally, it is not uncommon for files to be deleted from computing  
20 device 200. For example, the user may decide he or she no longer desires to run  
21 any programs that use a particular file (and uninstalls the program from the  
22 computing device), or the user no longer desires to keep a document file he or she  
23 created, etc. In these situations, a component of computing device 200 (e.g.,  
24 distributed system file interface 210) forwards an indication to one or more other  
25 computing devices 200 that the file has been deleted from computing device 200.

1 The other computing devices 200 that this indication is communicated to include  
2 the same computing devices that file information generation module 220  
3 previously determined the file identifier should be sent to, thereby allowing those  
4 devices to remove the file information entry from their respective databases.

5 The file information generated by a computing device is communicated to  
6 one or more computing devices referred to herein as database servers. Each  
7 database server maintains a database of file information that it receives and  
8 compares the received file information to identify any file information for two  
9 files that is the same (and thus indicative of potentially identical files). The  
10 database servers may be dedicated database servers (e.g., storing only file  
11 information), or alternatively may be other computing devices 200 in the network,  
12 storing both received file information as well as other files 240 in the distributed  
13 system portion(s) of their storage devices 208.

14 In a database server, file identification comparison module 242 receives file  
15 information and a corresponding file identifier (e.g., filename) from one or more  
16 other computing devices 200. Module 242 manages a database 244 (e.g., stored  
17 on device 208) of the file information it receives. Database 244 maintains a  
18 mapping of the file information to the file identifier. Database 244 may also  
19 maintain an indication of the computing device on which the file corresponding to  
20 the received file information is stored (or alternatively this may be inherent in the  
21 file identifier, which may include a filename as well as directory path to locate the  
22 file). Alternatively, the file identifier may not be stored (so long as the computer  
23 at which the file corresponding to the file information is stored is maintained in the  
24 database or otherwise known, the file information can be returned to that computer  
25 as an identification of the file). As discussed herein, the transferring of file

1 information from one computing device to another also typically encompasses  
2 transferring the file identifier as well.

3 Module 242 also compares the received file information to determine  
4 whether any of the previously received file information matches (e.g., whether two  
5 or more are the same). In one implementation, each time file information is  
6 received at the database server, module 242 compares the received file information  
7 to the database of file information 244 to determine whether a match exists.

8 If module 242 detects a file information match, then appropriate action is  
9 taken to move one or more of the files corresponding to the matching file  
10 information to the same computing device. Once the files corresponding to the  
11 matching file information are on the same computing device, the SIS component  
12 on that computing device is invoked to determine whether in fact the two files are  
13 identical, and if so then to delete one of the files and set up a pointer to the other  
14 file in its place. Module 242 can be responsible for moving files as necessary so  
15 that they are located on the same device, or alternatively this responsibility may be  
16 carried by the computing devices on which the potentially identical files are  
17 stored.

18 The copying of files to the same computer can be carried out in any of a  
19 wide variety of manners. In one implementation, module 242 forwards a  
20 command to one of the computers storing one of the files corresponding to the  
21 matching file information to relocate its file to the computer on which the other  
22 file corresponding to the matching file information is located. In another  
23 implementation, module 242 forwards the matching file information to the  
24 computing devices from which the matching file information were received, along  
25 with an indication that the match was identified. The individual computing

1 devices then coordinate with one another to transfer one of the files to the other  
2 computing device.

3 Fig. 3 is a flowchart illustrating an exemplary process followed to inform  
4 database servers of the file information in accordance with certain embodiments of  
5 the invention. The process of Fig. 3 is carried out by a computing device 200 of  
6 Fig. 2, and may be implemented in software.

7 Initially, the process waits until it is time to generate new file information  
8 for a file (act 250). Once it is time to generate new file information, the file for  
9 which the file information is to be generated is identified (act 252), and the file  
10 information is generated for that file (act 254). Optionally, the computing device  
11 may then store the generated file information and wait for additional file  
12 information to be generated (act 256), and then return to act 250 to generate more  
13 file information. The optional waiting period allows file information for multiple  
14 files to be forwarded to the identified database server(s) as a set rather than one-  
15 by-one. After the waiting period is over, or if the optional waiting is not  
16 performed, one or more database servers to receive the generated file information  
17 are identified (act 258). Which one or more database servers are to receive the file  
18 information can vary, as discussed in more detail below. The generated file  
19 information(s) is then transmitted to the identified database servers (act 260). It  
20 should also be noted that, based on different implementations as discussed below,  
21 the database server(s) to which the file information is to be transferred may not be  
22 readily identifiable (e.g., the computing device may not be aware of them yet).

23 Fig. 4 is a flowchart illustrating an exemplary process followed by a  
24 database server in accordance with certain embodiments of the invention. The  
25 process of Fig. 4 is carried out by a computing device 200 of Fig. 2, or

1 alternatively a dedicated server (e.g., a device 200 without file information  
2 generation module 220) and may be implemented in software.

3 Initially, file information is received (act 280). The manner in which the  
4 file information for various files is received (e.g., individually or in sets), as well  
5 as which computers the file information is received from, can vary and is  
6 discussed in more detail below. The received file information is optionally  
7 forwarded to one or more other database servers (act 282). Whether the file  
8 information is forwarded to another database server(s), as well as to what server(s)  
9 the file information is forwarded, varies by implementation as discussed in more  
10 detail below. Regardless of whether the file information is forwarded to other  
11 database servers, a check is made as to whether the file information should be  
12 added to the database of the database server that received the file information (act  
13 284). Whether the file information should be added to the database is based on  
14 certain criteria that vary by implementation, as discussed in more detail below. In  
15 some implementations, there is no checking in act 284 and all received file  
16 information is added to the database. If the received file information is not to be  
17 added to the database, then the process returns to act 280 where additional file  
18 information is eventually received.

19 However, if the file information is to be added to the database, then the  
20 received file information is added to the database maintained by the database  
21 server (act 286), and is compared to other file information in the database (act  
22 288). The database server also checks whether the newly received file information  
23 matches (is the same as) any of the file information already in the database (act  
24 290). If the received file information does match file information(s) in the  
25 database, then the computers storing the files corresponding to the matching file

1 information are notified of the match (act 292) so that they can take appropriate  
2 action. The process then returns to act 280 where additional file information is  
3 eventually received.

4 Additionally, in some situations copies of files may be replicated and stored  
5 in multiple locations (e.g., different computers) in the network for fault tolerance  
6 purposes. For example, in a serverless distributed file system, where the user has  
7 no guarantee that his or her file will be stored on a particular computer, the file  
8 may be replicated and stored on multiple computers so that the user can still access  
9 his or her file even if one or more of the computers is unavailable. When such  
10 replicated files exist in the network, care should be taken to ensure that they are  
11 not identified as duplicate copies and combined into a single file, and thus subvert  
12 the fault tolerance created by the replicated copies.

13 In one embodiment, the management of replicated file copies is handled by  
14 computing devices acting as directory servers (e.g., in distributed file system 150  
15 of Fig. 1). In this embodiment, the directory servers are responsible for both  
16 replicating files as well as identifying duplicate files, and thus know whether a  
17 particular file is a replica they created of another file. In one implementation, the  
18 duplicate identification is performed at a higher level than the replicated storage  
19 (e.g., duplicate identification is performed prior to replicating a file), thereby  
20 avoiding identification of a replicated file as a potentially duplicate file.

21 In certain embodiments discussed herein, various decisions are made by the  
22 computers based on a number of computers that exist in the network. Computers  
23 can determine an approximate number of computers that are coupled together in  
24 the network in a variety of conventional manners (note, however, that in some  
25 situations it is difficult to obtain an exact number of computers that are coupled

1 together in a network if the number of computers is very high, because computers  
2 can be continually joining and leaving the network). In one implementation, each  
3 time a computer logs into (or is otherwise coupled to) a network its presence is  
4 advertised to the network and propagated by the computers throughout the  
5 network. Additionally, each time a computer logs off (or is otherwise de-coupled  
6 from) a network, its retirement is advertised to the network and propagated by the  
7 computers throughout the network. Additional monitoring computers may also be  
8 established to monitor computers coupled to the network and detect (e.g., due to  
9 inactivity) their retirement from the network. Alternatively, any of a variety of  
10 other conventional processes may be used for identifying the topology and/or  
11 number of computers in the network.

12 Various different implementations for forwarding the file information to a  
13 database server(s), as well as communication among multiple database servers,  
14 exist. These various implementations will now be discussed. It should be noted  
15 that, in the discussions herein, reference is made to client computers and database  
16 server computers. These references are for the purposes of communicating and  
17 managing file information as described herein. In the distributed serverless  
18 environment, computers can be both client computers as well as database server  
19 computers.

### 21 **Centralized Database Implementation**

22 In the centralized database implementation, the client computers in the  
23 network are categorized into one or more groups, and each group includes one or  
24 more database servers. For each group, each client computer in that group  
25 forwards the file information it generates to one or more of the database servers in



1 that group. Each database server can then identify potentially identical files based  
2 on the file information it receives from client computers in that group.  
3 Additionally, the servers may optionally forward the file information they receive  
4 to other servers in other groups, thereby allowing potentially identical files located  
5 on client computers that have been categorized into different groups to be  
6 identified.

7 Fig. 5 illustrates an exemplary centralized database implementation in  
8 accordance with certain embodiments of the invention. In the illustrated example,  
9 a network 300 of multiple client computers (C) are categorized into multiple ( $n$ )  
10 groups 302, 304, and 306. Each group may include the same number of client  
11 computers (C), or alternatively different numbers. Furthermore, each client  
12 computer (C) belongs to one group, and may optionally belong to multiple groups  
13 (resulting in the client computer forwarding its file information to database servers  
14 for multiple groups).

15 Each group 302, 304, and 306 also includes one or more database servers  
16 (S). Although only one database server is illustrated in each group of Fig. 5,  
17 multiple database servers may be included in any one or more of the groups 302,  
18 304, and 306. Each group 302, 304, and 306 may include the same number of  
19 database servers, or alternatively varying numbers of database servers. The  
20 database servers in the groups 302, 304, and 306 communicate with each other,  
21 with each database server transferring the file information it receives to the  
22 database servers of the other groups. This communication among the database  
23 servers allows the file information to be shared, so that potentially identical files  
24 stored on client computers (C) in different groups can be identified.  
25

1 The manner in which client computers are categorized or separated into  
2 groups can vary. In one embodiment, the categorization is based on the naming  
3 convention used in naming the client computers and servers in network 300. The  
4 naming convention used in network 300 establishes multiple namespace roots  
5 which are assigned to selected client computers or servers in network 300, and  
6 then multiple lower-level names that are "under" the corresponding namespace  
7 root computers. One or more of these namespace root client computers or servers,  
8 as well as all of the lower-level names under those roots, belong to the same  
9 group.

10 Alternatively, client computers can be categorized into different groups in  
11 different manners, such as randomly, by client computer type, based on the date  
12 and/or time that they were coupled to network 300, based on geographic location,  
13 based on network connection type, and so forth.

14 Each client computer (C) knows the server (S) to which it is to transfer the  
15 file information it generates. In one embodiment, each client computer (C)  
16 transfers the file information it generates to the system at its namespace root,  
17 which is a database server (S). Alternatively, each client computer (C) may be  
18 programmed in another manner with an indication of the server (S) to which it is  
19 to transfer the file information it generates. The client computer (C) may receive a  
20 communication from a namespace server (S) identifying where the client computer  
21 should transfer its file information, or alternatively the client computer (C) may  
22 locate the database server (S) itself. For example, the group with the namespace  
23 root corresponding to a client computer (C) may keep information (e.g., addresses)  
24 identifying the database servers (S) for the group the computer (C) is in. The  
25 namespace root computer may identify all of these database servers (S) to the

1 requesting client computer (C), or alternatively may assign the client computer (C)  
2 to communicate with a particular one of the database servers (S). Additionally, a  
3 client computer (C) may communicate with one or more other client computers  
4 (C) to identify the root (or other) computer that it needs to access to determine the  
5 database server (S) to which it is to transfer the file information it generates.

6 When multiple database servers (S) exist within a group, responsibility for  
7 managing the database can be shared by the servers in any of a variety of manners.  
8 For example, particular servers may be assigned to receive file information from  
9 client computers (C) in particular address ranges, or file information for files in  
10 particular size ranges or creation date ranges, and so forth. This allows load and  
11 storage requirements to be partitioned among multiple database servers.

12 Additionally, multiple database servers may be employed for fault  
13 tolerance. In this situation, multiple servers are assigned to the same file  
14 information range so that if one or more of the servers fails (or is otherwise  
15 inaccessible) another is still available to do the processing. When employing  
16 multiple database servers for fault tolerance, care should be taken so that all of the  
17 servers handling a particular file information range are coordinated so as to  
18 generate only a single message to the client machines informing them about the  
19 detection of a potentially duplicate file. Alternatively, clients may only send file  
20 information to a single server assigned to a file information range and then rely on  
21 the servers to notify each other of new file information that any one of them has  
22 received. If a client cannot reach one server then it tries another assigned to the  
23 range.

24 Fig. 6 illustrates an exemplary network environment 300 in which  
25 responsibility for managing the file information database is shared among multiple

1 database servers in each group. Although one or more client computers (C) exist  
2 in each group 302, 304, and 306, for ease of explanation and to avoid cluttering  
3 the drawings the client computers (C) have not been shown. Rather, only the  
4 database servers (S) are illustrated in the groups 302, 304, and 306.

5 In the illustrated example of Fig. 6, each group 302, 304, and 306 includes  
6 the same number ( $k$ ) of database servers (S). Alternatively, each group need not  
7 include the same number of database servers (S). For example, a set of rules or an  
8 algorithm could be defined that tells each database server (S) in a group which one  
9 or more database servers (S) in the other groups to communicate with (e.g., group  
10 302 might have twice as many database servers (S) as group 304, with the file  
11 information space being divided up so that the piece that a database server (S) in  
12 group 304 handles is equivalent to two pieces handled by two different database  
13 servers (S) in group 302). By way of another example, if communication between  
14 database servers (S) in different groups is not needed, then each group need not  
15 include the same number of database servers (S). For purposes of discussion,  
16 however, it is assumed that each group 302, 304, and 306 includes the same  
17 number of database servers.

18 The file information generated by a client computer (C) is used to  
19 determine which database server (S) to transmit the file information to. After  
20 generating the file information, the client computer (C) calculates the following  
21 value:

$$22 \quad v = \text{info} \bmod k$$

23 where *info* is the generated file information and  $k$  is the number of database  
24 servers (S) in the group. The resultant value  $v$  is a value ranging from zero to  
25  $(k - 1)$ . Each of the  $k$  database servers is associated with one of the values in the

1 range from zero to  $(k-1)$ , and the client computer (C) forwards the file  
2 information to the database server associated with the resultant value  $v$ .

3 Each of the database servers (S) also communicates with the corresponding  
4 database servers (S) in the other groups. By identifying the database server (S)  
5 that is to handle particular file information based on the file information itself, the  
6 number of database servers (S) in the other groups that need to be communicated  
7 with in order to identify potential duplicate files across different groups is reduced  
8 (basically, each server need only communicate with one other server in each other  
9 group). So, for example, if a client computer (C) in group 304 generates file  
10 information that results in a value  $v$  of zero, the client computer communicates the  
11 file information to server 310. Server 310 is then able to compare the received file  
12 information to other file information it stores and identify any potential duplicate  
13 files within group 304. Additionally, database server 310 communicates the  
14 generated file information to servers 312 and 314 to identify any potential  
15 duplicate files in groups 306 and 302, respectively.

16 When database servers (S) communicate with database servers (S) in other  
17 groups, the file information sent between groups is not stored by the servers in the  
18 other groups (because it does not represent information about files in their groups).  
19 Rather, the file information is used to identify any matches with file information  
20 stored by the receiving database server, and then dropped after the match checking  
21 is completed. Alternatively, the file information could be stored by the servers in  
22 other groups (optionally with an indication of from which other group the file  
23 information was received).

24 In the centralized database implementation, two special situations arise.  
25 One situation is when the number of client computers in each group is equal to

one, and the other is when the number of client computers in a group is equal to all of the computers in the network. These special situations will now be discussed.

If the number of computers in each group is equal to one, then the centralized database implementation becomes a "localized pair-wise checking" implementation in which each of the client computers is its own group and each client computer maintains its own file information mappings. Thus, each client computer also acts as a database server. Whenever one client computer becomes aware of another client computer in the network, the client computer communicates all of the file information it has generated for its files to the other client computer, allowing the other client computer to check for potentially duplicate files. The communication may occur immediately after the client computer becomes aware of the other client computer, or alternatively after a period of time (e.g., a delay may be incurred while the computer is performing other functions, while the computer waits for a period of low use on the network, and so forth). The client also subsequently sends incremental file information updates to the other client computer as new file information is generated.

Fig. 7 illustrates the localized pair-wise checking implementation in additional detail. For ease of explanation, network 350 is illustrated including only nine computers (C). In network 350, computers  $C_1$  and  $C_2$  are aware of each other and have communicated their file information between them. Similarly, computers  $C_2$  and  $C_3$  are aware of each other and have communicated their file information between them. Note, however, that the computers  $C_1$  and  $C_3$  are not aware of each other and thus have not communicated their file information between them. Additionally, computers  $C_1$  and  $C_4$  are aware of each other, as are computers  $C_1$  and  $C_5$ , and computers  $C_4$  and  $C_5$ . Thus, for computers  $C_1$  through

1 C<sub>5</sub>, each of the computers is aware of some of the other computers C<sub>1</sub> through C<sub>5</sub>,  
2 but not all.

3 For computers C<sub>6</sub>, C<sub>7</sub>, C<sub>8</sub>, and C<sub>9</sub>, each of these four computers is aware of  
4 each of the others, and thus each has communicated its file information to the  
5 others. Note, however, that none of the computers C<sub>1</sub> through C<sub>5</sub> is aware of any  
6 of the computers C<sub>6</sub> through C<sub>9</sub>, nor are any of the computers C<sub>6</sub> through C<sub>9</sub> aware  
7 of any of the computers C<sub>1</sub> through C<sub>5</sub>.

8 A client computer can become aware of another client computer in any of a  
9 wide variety of conventional manners. In one implementation, any of a variety of  
10 well-known network mapping processes can be used by a client computer to  
11 identify other client computers on the network it is coupled to. Alternatively, a  
12 computer may broadcast its presence when added to a network.

13 Alternatively, rather than forwarding its file information to any other client  
14 computer that a client computer becomes aware of, additional restrictions on what  
15 client computers the file information will be forwarded to may be imposed. For  
16 example, a client computer may forward its file information only to client  
17 computers that are within a particular range (e.g., geographically close, within a  
18 particular number of links or routers on the network, and so on).

19 In addition to transmitting its own file information to other client computers  
20 of which a particular client computer is aware, the client computer may also  
21 forward file information that it has received from other computers as well. For  
22 example, in network 350, client computer C<sub>2</sub> may initially become aware of client  
23 computer C<sub>3</sub>, and receive all of the file information of client computer C<sub>3</sub>. When  
24 client computer C<sub>2</sub> subsequently becomes aware of client computer C<sub>1</sub>, client  
25

1 computer C<sub>2</sub> communicates all of its file information, as well as all of the file  
2 information received from client computer C<sub>3</sub>, to client computer C<sub>1</sub>.

3 In one implementation, file information for each file is also associated with  
4 a "time to live" component that identifies how many client computers the file  
5 information can be communicated to. Each time the file information is  
6 communicated to another client computer, the time to live component is  
7 decremented by one. Once the time to live component reaches zero, the file  
8 information is not communicated to any more client computers. Various  
9 alternatives may be implemented for the time to live component, such as different  
10 threshold values could be used for different computers or different files (e.g., a file  
11 with an indicated or perceived greater importance could be assigned a larger value  
12 for its time to live component), the count could be decremented by more or less  
13 than one, the count could be incremented and compared to an upper bound rather  
14 than decremented and compared to zero, and so forth. For example, following the  
15 previous example, assume that the file information for each file from client  
16 computer C<sub>3</sub> has a time to live component with a value of two. When the file  
17 information is communicated to client computer C<sub>2</sub> the associated time to live  
18 component(s) for the file information of client computer C<sub>3</sub> on client computer C<sub>2</sub>  
19 are decremented to the value of one. Then, when the file information of client  
20 computer C<sub>3</sub> are communicated to client computer C<sub>1</sub> the associated time to live  
21 component(s) for the file information of client computer C<sub>3</sub> on client computer C<sub>1</sub>  
22 are decremented to the value of zero. Thus, even though client computer C<sub>1</sub> may  
23 be aware of, or may subsequently become aware of, client computers C<sub>4</sub> and C<sub>5</sub>,  
24 client computer C<sub>1</sub> does not communicate the file information of client computer  
25 C<sub>3</sub> to either of computers C<sub>4</sub> or C<sub>5</sub>. However, if client computer C<sub>3</sub> were to



1 subsequently become aware of either client computer  $C_4$  or  $C_5$ , then client  
2 computer  $C_3$  would communicate its file information to the appropriate one of  
3 client computer  $C_4$  and  $C_5$  and the associated time to live component(s) for the file  
4 information of client computer  $C_3$  on client computer  $C_4$  or  $C_5$  would be  
5 decremented to the value of one.

6 File information for each file may be associated with its own "personal"  
7 time to live component, or alternatively file information for multiple files from the  
8 same client computer may be grouped together (e.g., into a single set for the client  
9 computer) and have an associated time to live component. File information for  
10 different files and/or different computers can optionally have different time to live  
11 components. For example, file information for larger files may have longer time  
12 to live components than shorter files (e.g., assuming that the potential space  
13 savings of finding a duplicate of the larger file is worth the extra burden of  
14 communicating the file information to additional client computers).

15 Additionally, in the localized pair-wise checking implementation, file  
16 information can optionally be communicated among the computers in a  
17 compressed form. Any of a variety of conventional techniques can be used to  
18 communicate the information in a compressed form, such as the use of well-  
19 known Bloom filters. For additional information on Bloom filters, the reader is  
20 directed to L. Fan, P. Cao, J. Almeida, and A. Broder, "Summary Cache: A  
21 Scalable Wide-Area Web Cache Sharing Protocol", ACM SIGCOMM, 1998.

22 The other special situation that can arise in the centralized database  
23 implementation is when the number of client computers in a group is equal to all  
24 of the computers in the network. In this situation, the centralized database  
25 implementation reduces to a single group and the one or more database servers in

1 the network receive the file information from all the client computers. Each  
2 database server may receive file information from all of the computers, or  
3 alternatively only for select client computers (e.g., based on the file signature itself  
4 analogous to the discussion above regarding Fig. 6).

5 Fig. 8 illustrates the special situation when the number of client computers  
6 in a group is equal to all of the computers in the network in additional detail. In  
7 network 360, multiple (a) client computers (C) are illustrated along with multiple  
8 (b) servers (S). All of the client computers (C) are part of the same group,  
9 communicating their file information to one or more of the servers (S).

#### 11 **Fully Distributed Stochastically Partitioned Database Implementation**

12 In the fully distributed stochastically partitioned database implementation,  
13 each computer in the network operates as both a client computer and a database  
14 server. Alternatively, some machines might function only as clients and not as  
15 database servers, while other might function only as database servers and not as  
16 clients. Each computer generates file information for files stored at its computer,  
17 and forwards that generated file information to one or more other computers. To  
18 which computers particular file information is forwarded is based on both the  
19 generated file information as well as identifiers (ID's) for each computer in the  
20 network, as discussed in more detail below. Each computer, then, is responsible  
21 for comparing the file information it receives from computers in the network and  
22 determining whether any of the received file information matches each other.

23 In the fully distributed stochastically partitioned database implementation,  
24 each computer in the network is assigned a computer ID. The computer ID's can  
25 be assigned in any of a variety of manners. However, in order to spread out the

1 file information processing relatively evenly among all of the computers, the  
2 computer ID's should be assigned such that the computer ID's are fairly evenly  
3 distributed throughout the Hamming space of possible computer ID's. More  
4 specifically, this even distribution is important for a particular subset of  $W$  bits of  
5 the computer ID, as described below.

6 In one embodiment, each computer in the network includes a public/private  
7 key pair used in public key cryptography. The computer ID for a particular  
8 computer is generated based on the public key of this key pair, such as by applying  
9 a one-way hashing function (e.g., SHA, MD5, etc.) to the public key and using the  
10 resultant hash value as the computer ID. Alternatively, different processes can be  
11 used to create the computer ID for a computer, such as use of a conventional  
12 random number generator (or pseudo-random number generator) by a central  
13 authority that assigns computer ID's, use of an identification number assigned to  
14 the CPU in the computer, and so forth.

15 Fig. 9 is a flowchart illustrating an exemplary process followed by each  
16 computer for the fully distributed stochastically partitioned database  
17 implementation in accordance with certain embodiments of the invention. The  
18 process of Fig. 9 is carried out by a computing device 200 of Fig. 2, and may be  
19 implemented in software.

20 For each file stored at the computer for which file information is generated,  
21 an imprint for the file is identified using  $W$  bits of the file information (act 380).  
22 Which  $W$  bits of the file information to use can vary, but should be consistent  
23 across all the files in the system. In one implementation, the  $W$  least significant  
24 bits of the file information are used as the imprint. The choice of which  $W$  bits to  
25

1 use should try to result in a fairly uniform mapping of imprint to file information  
2 so that unwanted clustering effects do not arise.

3 The computer also identifies each known computer in the network that has  
4 a computer ID that has the same  $W$  bits as the imprint (act 382). Which  $W$  bits of  
5 the computer ID to use can vary, but should be consistent across all the computers  
6 in the network. In one implementation, the  $W$  least significant bits of the computer  
7 ID are used. The choice of which  $W$  bits to use should try to result in a fairly  
8 uniform mapping of  $W$  bits to computer ID so that unwanted clustering effects do  
9 not arise. Alternatively other bits may be used (the selected bits of the file  
10 information used may be the same as the bits used for the imprint of the file  
11 information, or alternatively different bits may be selected). Once these  
12 computers are identified, the computer that generated the file information sends  
13 the file information to each of the computers identified in act 382 (act 384).

14 Each computer calculates its own value of  $W$  as follows:

$$15 \quad W = \left\lfloor \lg \frac{M}{R} \right\rfloor$$

16 where the value  $M$  is the total number of computers in the network that the  
17 computer knows about (possibly including itself),  $R$  is a system configuration  
18 parameter,  $\lg$  indicates a binary (base 2) logarithm, and the floor brackets indicate  
19 the largest integer that is no greater than the enclosed value. The value  $M$   
20 represents the number of computers that function as database servers; if some  
21 machines act solely as clients and not as database servers, then they will not be  
22 included in this number. The value  $M$  can vary by computer, which means that the  
23 value  $W$  can vary by computer. However, despite these variations, potentially  
24 identical files can still be identified. Each computer can identify the value  $M$  in  
25

1 any of a wide variety of conventional manners, such as using any of a variety of  
2 conventional network topology identification processes to determine the location  
3 and number of computers in the network.

4 The value  $R$  is a system configuration parameter that imposes a bound on  
5 the average number of computers to which particular file information is  
6 communicated. The bound imposed by  $R$  is as follows:

$$7 \quad R \leq \lambda < 2R$$

8 where  $\lambda$  is the average number of computers to which particular file information  
9 is communicated. The value of  $R$  can vary by implementation. In one  
10 implementation, typical values for  $R$  range from 3 to 6.

11 Fig. 10 illustrates an exemplary network 400 in which the fully distributed  
12 stochastically partitioned database implementation is employed. Although  
13 network 400 includes many computers, only five computers are illustrated in Fig.  
14 10 for ease of explanation and to avoid cluttering the drawings. Network 400  
15 includes computers 402, 404, 406, 408, and 410. The communication of file  
16 information for two files from each of computers 402 and 404 is illustrated in Fig.  
17 10.

18 In the example of Fig. 10, assume that computers 402 and 404 each believe  
19 a different number of computers exist in network 400, and that computer 402 has  
20 calculated a value of  $W=2$ , while computer 404 has calculated a value of  $W=3$ .  
21 Further assume that the location of the  $W$  bits being used for both the file  
22 information and the computer ID's are the  $W$  least significant bits. Each of the  
23 computers 402 – 410 is assigned a computer ID. Only the three least significant  
24 bits of the computer ID is shown for each computer 402 – 410; the more  
25 significant bits of the computer ID are not shown. As illustrated, the least

1 significant bits of the computer ID for computers 402 and 406 are "000", while the  
2 least significant bits of the computer ID for computer 404 are "010", the least  
3 significant bits of the computer ID for computer 408 are "100", and the least  
4 significant bits of the computer ID for computer 410 are "101".

5 Two files 412 and 414 are illustrated as stored at computer 404, having file  
6 information with least significant bits of "000" and "100", respectively. Computer  
7 404 has calculated a value of  $W=3$ , so computer 404 generates an imprint for file  
8 412 that is the three least significant bits of the file information for file 412. The  
9 imprint of file 412 is thus "000". Computer 404 then transfers the file information  
10 for file 412 to all other computers in network 400 that have a computer ID with the  
11 three least significant bits equal to "000". Thus, computer 404 transfers the file  
12 information for file 412 to computer 402 and computer 406. Similarly, the imprint  
13 of file 414 is "100", so computer 404 transfers the file information for file 414 to  
14 computer 408.

15 Two additional files 416 and 418 are illustrated as stored at computer 402,  
16 having file information with least significant bits of "100" and "000", respectively.  
17 Computer 402 has calculated a value of  $W=2$ , so computer 402 generates an  
18 imprint for file 416 that is the two least significant bits of the file information for  
19 file 416. The imprint of file 416 is thus "00". Computer 402 then transfers the file  
20 information for file 416 to all other computers in network 400 that have a  
21 computer ID with the two least significant bits equal to "00". Thus, computer 402  
22 transfers the file information for file 416 to computer 406, computer 408, and  
23 computer 402 (back to itself). Similarly, the imprint of file 418 is also "00", so  
24 computer 402 also transfers the file information for file 418 to computers 406,  
25 408, and 402.

1 It should be noted that in the example of Fig. 10, computers 402 and 404  
2 have calculated different values of  $W$ . This results in computers 402 and 404  
3 identifying different imprints for their file information and sending them to  
4 different sets of computers (e.g., even though the least significant bits of both files  
5 412 and 418 are "000", the file information for file 412 (having an imprint of  
6 "000") is not sent to computer 408, while the file information for file 418 (having  
7 an imprint of "00") is sent to computer 408). Essentially, computer 402 ends up  
8 typically sending its file information to more computers than computer 404.  
9 However, potentially identical files on computers 404 and 402 can still be  
10 identified because the set of computers derived from a smaller value of  $W$  is a  
11 superset of those derived from a larger value of  $W$  (so the file information from  
12 both computers is sent to some of the same computers (e.g., computers 402 and  
13 406)).

14 It should also be noted that situations can arise where there is no computer  
15 with a computer ID that has the  $W$  bits matching the imprint of the file  
16 information. For example, if  $W=3$ , and the imprint is "001", situations can arise  
17 where there are no computers having a computer ID with the corresponding bit  
18 values of "001". In one implementation, this situation is resolved by simply not  
19 forwarding the file information to any computer. However, note that in the  
20 example of Fig. 10, computer 402 has calculated a value of  $W=2$ , so it would send  
21 any file information that ends with "001" to computer 410, since the two least  
22 significant bits match. Thus, although calculating a lower value of  $W$  increases the  
23 work that a computer does (as described above), it also increases the probability  
24 that duplicate files will be found. Alternatively, other solutions may be used when  
25 there is no identified computer for some values of file information, such as

1 assigning a particular computer to be the recipient of any such file information, or  
2 changing one or more bits of the imprint (so long as all the computers agree to use  
3 the same algorithm for changing the bits of the imprint).

#### 4 5 **Group-Based System Using Directory Services Implementation**

6 The group-based system using directory services implementation is similar  
7 to the fully distributed stochastically partitioned database implementation.  
8 Imprints are generated based on file information as discussed above, however, a  
9 database of imprint to computer ID mappings is accessed to determine which  
10 computers the file information is to be communicated to, thereby requiring the file  
11 information to potentially be sent to fewer computers than in the fully distributed  
12 stochastically partitioned database implementation.

13 Fig. 11 is a flowchart illustrating an exemplary process followed by each  
14 computer for the group-based system using directory services implementation in  
15 accordance with certain embodiments of the invention. The process of Fig. 11 is  
16 carried out by a computing device 200 of Fig. 2, and may be implemented in  
17 software.

18 For each file stored at the computer for which file information is generated,  
19 an imprint for the file is identified using  $W$  bits of the file information (act 440),  
20 analogous to act 380 of Fig. 9 above. An imprint to computer mapping is then  
21 accessed (act 442). The imprint to computer mapping is initially retrieved from  
22 one or more computers in the network that are designated mapping servers. The  
23 mapping may optionally be subsequently cached at the computer so that  
24 subsequent requests can be handled by the computer locally rather than requiring a  
25 network access. Based on this mapping, one or more computers in the network to



1 which the file information is to be transferred are identified (act 444), and the  
2 computer sends the file information to those other computers (act 446). The  
3 imprint to computer mapping may map the imprint to a computer ID, or  
4 alternatively some other name or identification of the computer.

5 The imprint to computer ID mappings are stored on the designated mapping  
6 servers and are accessible to other computers in the network. The designated  
7 mapping servers may be dedicated mapping servers, or alternatively may be  
8 computing devices such as device 200 of Fig. 2 that include both server and client  
9 functionality. Analogous to the database servers discussed above, multiple  
10 computers may be designated mapping servers, and each computer knows one or  
11 more mapping servers (or can ascertain the identity of one or more mapping  
12 servers) from which it can retrieve mappings. Also analogous to the database  
13 servers discussed above, if multiple designated mapping servers are employed,  
14 they may share mapping information (for fault tolerance purposes, such as one  
15 being a backup for another), or alternatively different servers may be designated to  
16 handle requests for different imprints (for load sharing purposes).

17 The imprint to computer mapping maps the imprint to one or more  
18 computers in the network. The imprint to computer mapping may map the imprint  
19 to a computer(s) having a computer ID that has the same  $W$  bits as the imprint, or  
20 alternatively a computer having a computer ID with  $W$  bits that are not the same as  
21 the imprint. In other words, there may be, but need not be, any correlation  
22 between the imprint and the  $W$  bits of the computer ID's in the mapping. By not  
23 tying the mapping to the  $W$  bits of the computer ID, the mapping server(s) need  
24 not store information about the  $W$  bits of all computer ID's in the network. Rather,  
25

1 the mapping server(s) can store only the computer ID's of the set of computers that  
2 they have designated to be file information processing servers..

3 In one implementation, the computer sends the file information to each  
4 other computer identified in the imprint to computer mappings (act 446 of Fig.  
5 11). Alternatively, the computer may send the file information to only one of the  
6 computers identified in the imprint to computer mappings. According to this  
7 alternative, computers that receive the file information know which other  
8 computers are responsible for checking for file information matches for particular  
9 imprints (e.g., by accessing a designated mapping server and obtaining the imprint  
10 to computer mappings for that mapping). Any file information received by one of  
11 the computers is then forwarded to the other computer(s) responsible for checking  
12 for file information matches for that particular imprint.

13 Fig. 12 illustrates an exemplary network 460 in which the group-based  
14 system using directory services implementation is employed. Although network  
15 460 includes many computers, only five computers are illustrated in Fig. 12 for  
16 ease of explanation and to avoid cluttering the drawings. Network 460 includes  
17 computers 462, 464, 466, 468, and 470. The communication of file information  
18 for two files 472 and 474 from computer 470 is illustrated in Fig. 12.

19 In the example of Fig. 12, assume that computer 470 has calculated a value  
20 of  $W=3$ , and that the location of the  $W$  bits being used for both the file information  
21 and the computer ID's are the  $W$  least significant bits. Each of the computers 462  
22 – 470 is assigned a computer ID. Only the three least significant bits of the  
23 computer ID is shown for each computer 462 – 470; the more significant bits of  
24 the computer ID are not shown. As illustrated, the least significant bits of the  
25 computer ID for computers 462, 464, 466, and 470 are "000", while the least

1 significant bits of the computer ID for computer 468 is "010". Computer 468 is  
2 designated as the mapping server.

3 When computer 470 generates the file information for file 472, it uses the  
4 *W* least significant bits of the file information as the imprint, which is "000".  
5 Assuming computer 470 does not have a locally stored computer mapping for  
6 imprint "000", computer 470 sends a request 476 to mapping server 468 for the  
7 imprint to computer mapping for imprint "000". The mapping 478 is returned by  
8 mapping server 468, and stored in mappings 480 of computer 470. All computers  
9 identified by mapping 478 may be stored in mappings 480, or alternatively only a  
10 subset of the computers (e.g., one or two computers). For purposes of discussion,  
11 assume that mapping 478 indicates that computers 466 and 462 are to receive file  
12 information with imprints of "000". Computer 470 then forwards the file  
13 information 482 for the file 472 to computer 466, which in turn receives the file  
14 information 482 and communicates it to computer 462. Alternatively, computer  
15 470 may forward the file information 482 to both computer 462 and 466.

16 Subsequently, computer 470 generates the file information for file 474 and  
17 identifies the imprint of the file information as "000". Rather than accessing  
18 mapping server 468, local mapping 480 is accessed to identify that the file  
19 information is to be communicated to computer 466 (and/or computer 462).  
20 computer 470 then forwards the file information 484 to computer 466 (and/or  
21 computer 462). If the computer identified in mapping 480 is not available (e.g.,  
22 computer 462 is identified in mapping 480 but it has failed or is otherwise  
23 inaccessible), computer 470 sends another request to mapping server 468  
24 requesting identification of another computer(s) that is mapped to the imprint  
25 "000".

1 It should be noted that, by having more than one computer in the mapping  
2 for a given  $W$ , improved fault tolerance is achieved because there are multiple  
3 computers that can process the file information. However, if fault tolerance is less  
4 of a concern, then fewer (including only one) computers may be included in the  
5 mapping for a given  $W$ .

### 6 7 **Multi-Level Stochastically Partitioned Database Implementation**

8 The multi-level stochastically partitioned database implementation is  
9 similar to the fully distributed stochastically partitioned database implementation.  
10 Imprints are generated based on file information as discussed above, however,  
11 similar to the group-based system using directory services implementation, the  
12 number of computers to which the file information are communicated to is  
13 reduced.

14 The multi-level stochastically partitioned database implementation can be  
15 employed using any number of levels, and is described herein primarily with  
16 reference to a two-level implementation. In a two-level implementation, the  
17 number of other computers that each computer has to contact to forward its file  
18 information to is proportional to the square root of the system size, while in a  
19 three-level implementation the number of computers that each computer has to  
20 contact to forward its file information to is proportional to the cube root of the  
21 system size. Alternative embodiments may also be used, with the number of  
22 computers that each computer has to contact to forward its file information to  
23 grows non-polynomially proportionally to the system size (e.g., based on  
24 logarithmic growth).

Fig. 13 is a flowchart illustrating an exemplary process followed by each computer for the multi-level stochastically partitioned database implementation in accordance with certain embodiments of the invention. The process of Fig. 13 is carried out by a computing device 200 of Fig. 2, and may be implemented in software. The process of Fig. 13 is carried out by each computer in the network, and is discussed with reference to a current computer (the computer, as discussed, that is determining to what computer to forward file information).

Initially, a value of  $W$  is identified based on the number of computers in the network (act 500), analogous to the discussions above regarding the fully distributed stochastically partitioned database implementation. A first group of computers, referred to as the group  $S_B$  is identified as the computers having the same  $W$  bits as the current computer ID (act 502). This group  $S_B$  thus includes the current computer. A second group of computers, referred to as the group  $S_0$  is identified as the computers having a first subset of the  $W$  bits the same as the current computer ID (act 504). In one implementation, the first subset of the  $W$  bits are the even bits of the  $W$  bits. This first subset can vary (e.g., it may be the odd bits, or in a three-level implementation two out of every three bits starting with bit zero, etc.). A third group of computers, referred to as the group  $S_1$  is identified as the computers having a second subset of the  $W$  bits the same as the current computer ID (act 506). This second subset can also vary, but is different than the first subset (e.g., it may be the even bits, or in a three-level implementation two out of every three bits starting with bit one, etc.). Although only three groups are illustrated as being identified in Fig. 13, additional groups are created for each additional level in the implementation, and the computers in those groups identified in an analogous manner. For example, in a three-level

1 implementation a fourth level is identified based on a third subset of the  $W$  bits  
2 (e.g., the computers having two out of every three bits of the  $W$  bits starting with  
3 bit two the same as the current computer ID).

4 These identified groups are then used in determining, for each file  
5 information being analyzed by the current computer, to which computers to send  
6 the file information. Each file information being analyzed by the current computer  
7 may have been generated at the current computer, or alternatively may have been  
8 generated at another computer and transferred to the current computer. The  
9 computer then waits for new file information that needs to be analyzed (act 508).  
10 The groups identified in acts 502 – 506 may take on new members as new  
11 machines are added to the system, or they may lose members as machines are  
12 removed from the system, but otherwise they remain the same until the number of  
13 computers in the network changes sufficiently to alter the value of  $W$ . When this  
14 occurs, acts 502 – 506 are repeated to re-identify the groups.

15 Eventually, new file information is received, and an imprint for the file is  
16 identified (act 510), analogous to the discussions above regarding the fully  
17 distributed stochastically partitioned database implementation. The current  
18 computer then checks whether all the bits of the imprint match (are the same as)  
19 the  $W$  bits of the current computer ID (act 512). If all the bits do match, then the  
20 file information is transferred to each computer in the first group, group  $S_B$ , (act  
21 514). However, if all the bits do not match, then a check is made as to whether the  
22 second subset of bits of the imprint match (are the same as) the second subset of  
23 bits of the current computer ID (act 516). If the second subsets do match then the  
24 file information is transferred to each computer in the third group, group  $S_1$ ,  
25 having computer ID's with their  $W$  bits matching (the same as) the imprint (act

518). However, if the second subsets do not match, then the file information is transferred to each computer in the second group, group  $S_0$ , having computer ID's with the second subset of their  $W$  bits matching (the same as) the imprint (act 520).

Although the decision of to which group of computers to send the file information is based on two subsets of bits in acts 512 – 520, alternatively an additional subset of bits is analyzed for each additional level in the implementation. For example, in a three-level implementation a third subset of bits is analyzed if the second subsets of bits of the imprint do not match the second subset of bits of the current computer ID in act 516. Based on this third set of bits, then, the file information is either sent to selected computers in the fourth group, or to selected computers in the second group.

An exemplary process carried out by the current computer in determining to which computer(s) to forward its file information is described in more detail as follows. Initially, the value of  $W$  is calculated as discussed above:

$$W = \left\lceil \lg \frac{M}{R} \right\rceil$$

Two additional values,  $W_0$  and  $W_1$  are then calculated based on  $W$  as follows:

$$W_0 = \left\lceil \frac{W}{2} \right\rceil \quad W_1 = \left\lfloor \frac{W}{2} \right\rfloor$$

The ceiling brackets indicate that  $W_0$  is set to the smallest integer that is no less than  $W/2$ , and the floor brackets indicate that  $W_1$  is set to the largest integer that is no greater than  $W/2$ . The current computer then calculates two bit masks:  $\psi_0$  which is a sequence of  $W_0$  copies of the bit string "01", and  $\psi_1$  which is a sequence of  $W_1$  copies of the bit string "10". These can be illustrated as follows:

$$\psi_0 = \sum_{k=0}^{W_0-1} 2^{2k} \quad \psi_1 = \sum_{k=0}^{W_1-1} 2^{2k+1}$$

When the current computer becomes aware of a new computer in the network, the current computer may or may not remember the new computer's ID. If the current computer does decide to remember the new computer's ID, it assigns the new computer into a particular group. Whether to remember the new computer's ID, as well as to which group to assign the new computer, is determined according to the following pseudocode (where "&" refers to bitwise conjunction, "==" refers to "is equal to",  $CID_{new}$  refers to the  $W$  bits of the computer ID of the new computer, and  $CID_{me}$  refers to the  $W$  bits of the computer ID of the current computer).

```

if ( $CID_{new} \& \psi_0$ ) == ( $CID_{me} \& \psi_0$ ) and ( $CID_{new} \& \psi_1$ ) == ( $CID_{me} \& \psi_1$ )
    add new computer to set  $S_B$ ;
else if ( $CID_{new} \& \psi_0$ ) == ( $CID_{me} \& \psi_0$ )
    add new computer to set  $S_0$ ;
else if ( $CID_{new} \& \psi_1$ ) == ( $CID_{me} \& \psi_1$ )
    add new computer to set  $S_1$ ;
else
    forget new computer;

```

When new file information is identified (based on either a file stored at the current computer or file information received from another computer in the network), the current computer determines what to do with the file information based on the following pseudocode. In the following pseudocode, "==" refers to "is equal to", "!=" refers to "is not equal to", "&" refers to bitwise conjunction, "information" refers to the new file information,  $CID_{recipient}$  refers to the  $W$  bits of the computer ID of a computer that is a potential recipient of the new file information, and  $CID_{me}$  refers to the  $W$  bits of the computer ID of the current computer.



```

1   if (information &  $\psi_1$ ) != (CIDme &  $\psi_1$ )
2       send information to every computer in S0 for which
3       (information &  $\psi_1$ ) == (CIDrecipient &  $\psi_1$ );
4   else if (information &  $\psi_0$ ) != (CIDme &  $\psi_0$ )
5       send information to every computer in S1 for which
6       (information &  $\psi_0$ ) == (CIDrecipient &  $\psi_0$ );
7   else {
8       if information originated from current computer
9           send information to every computer in SB;
10          store information in database of current computer;
11          check database for matching information;
12          notify pairs of computers with matching information;
13      }
14

```

Fig. 14 illustrates an exemplary network 530 in which the multi-level stochastically partitioned database implementation is employed. The example network 530 includes only 32 computers for ease of explanation and to avoid cluttering the drawings. Additionally, only the five least significant bits of the computer ID (CID) for each computer is illustrated in Fig. 14.

Fig. 14 illustrates a two-level stochastically partitioned database described from the point of view of computer CID 11001. Assume, for the purposes of discussion of Fig. 14, that  $R=2$  and the following values have been computed:  $W=4$ ,  $W_0=2$ ,  $W_1=2$ ,  $\psi_0=0101$ , and  $\psi_1=1010$ . Based on these values, and the computer ID's illustrated in Fig. 14, computer CID 11001 groups selected machines into three groups as follows. Group S<sub>B</sub> includes computer CID 01001. Group S<sub>0</sub> includes the following computers: CID 00001, CID 00011, CID 01011, CID 10001, CID 10011, and CID 11011. Group S<sub>1</sub> includes the following computers: CID 01000, CID 01100, CID 01101, CID 11000, CID 11100, and CID 11101.

1 When new file information is identified, computer CID 11001 identifies the  
2  $W$  (4 in this example) least significant bits of the file information. If the four least  
3 significant bits of the file information are "1001", then the file information is  
4 stored in the database of computer CID 11001. The file information is also  
5 forwarded to other computers in group  $S_B$  (computer CID 01001), which also store  
6 the file information in their databases. The transfers to computers in group  $S_B$  are  
7 referred to as "zero-hop" transfers, and are illustrated by the dashed line from  
8 computer CID 11001 to computer CID 01001.

9 If the four least significant bits of the file information are " $1x0y$ " for any  
10 single-bit values of  $x$  and  $y$  other than  $(x,y) = (0,1)$ , then the file information is sent  
11 to computers in group  $S_1$  having CID's that are " $01x0y$ " or " $11x0y$ ". Upon receipt  
12 of the file information, these computers in group  $S_1$  will have the same four least  
13 significant bits of their CIDs matching the four least significant bits of the file  
14 information, so these computers will store the received file information in their  
15 respective databases. The transfers to computers in group  $S_1$  are referred to as  
16 "single-hop" transfers, and are illustrated by the single solid lines from computer  
17 CID 11001 to the computers in  $S_1$ .

18 If the four least significant bits of the file information are " $wxyz$ " for any  
19 single-bit values of  $w$ ,  $x$ ,  $y$  and  $z$  other than  $(w,y) = (1,0)$ , then the file information  
20 is sent to computers in group  $S_0$  having CID's that are " $0w0y1$ " and " $1w0y1$ ".  
21 Upon receipt of the file information, these computers in group  $S_0$  will either store  
22 the file information in their respective databases, or forward the file information  
23 on to another computer. If  $(x,z) = (0,1)$  then the four least significant bits of  
24 computers having CIDs " $0w0y1$ " and " $1w0y1$ " will match the four least significant  
25 bits of the file information, so these computers will store the file information in

1 their respective databases. However, if  $(x,z) \neq (0,1)$  then the computers with CIDs  
2 "0w0y1" and "1w0y1" will forward the file information to computers "0wxyz" and  
3 "1wxyz", which will in turn store the file information in their respective databases.  
4 The transfer to computers in group  $S_0$  are referred to as "double-hop" transfers  
5 because they may require a second transfer before reaching an appropriate  
6 database. These transfers are illustrated by the double solid lines from computer  
7 CID 11001 to the computers in  $S_0$ .

### 8 9 Example Computer System

10 Fig. 15 illustrates a more general exemplary computer environment 600,  
11 which can be used in various embodiments of the invention. The computer  
12 environment 600 is only one example of a computing environment and is not  
13 intended to suggest any limitation as to the scope of use or functionality of the  
14 computer and network architectures. Neither should the computer environment  
15 600 be interpreted as having any dependency or requirement relating to any one or  
16 combination of components illustrated in the exemplary computer environment  
17 600.

18 Computer environment 600 includes a general-purpose computing device in  
19 the form of a computer 602. Computer 602 can be, for example, any of computing  
20 devices 102-108 of Fig. 1, or a computing device 200 of Fig. 2. The components  
21 of computer 602 can include, but are not limited to, one or more processors or  
22 processing units 604, a system memory 606, and a system bus 608 that couples  
23 various system components including the processor 604 to the system memory  
24 606.

1 The system bus 608 represents one or more of any of several types of bus  
2 structures, including a memory bus or memory controller, a peripheral bus, an  
3 accelerated graphics port, and a processor or local bus using any of a variety of  
4 bus architectures. By way of example, such architectures can include an Industry  
5 Standard Architecture (ISA) bus, a Micro Channel Architecture (MCA) bus, an  
6 Enhanced ISA (EISA) bus, a Video Electronics Standards Association (VESA)  
7 local bus, and a Peripheral Component Interconnects (PCI) bus also known as a  
8 Mezzanine bus.

9 Computer 602 typically includes a variety of computer readable media.  
10 Such media can be any available media that is accessible by computer 602 and  
11 includes both volatile and non-volatile media, removable and non-removable  
12 media.

13 The system memory 606 includes computer readable media in the form of  
14 volatile memory, such as random access memory (RAM) 610, and/or non-volatile  
15 memory, such as read only memory (ROM) 612. A basic input/output system  
16 (BIOS) 614, containing the basic routines that help to transfer information  
17 between elements within computer 602, such as during start-up, is stored in ROM  
18 612. RAM 610 typically contains data and/or program modules that are  
19 immediately accessible to and/or presently operated on by the processing unit 604.

20 Computer 602 may also include other removable/non-removable,  
21 volatile/non-volatile computer storage media. By way of example, Fig. 15  
22 illustrates a hard disk drive 616 for reading from and writing to a non-removable,  
23 non-volatile magnetic media (not shown), a magnetic disk drive 618 for reading  
24 from and writing to a removable, non-volatile magnetic disk 620 (e.g., a "floppy  
25 disk"), and an optical disk drive 622 for reading from and/or writing to a

1 removable, non-volatile optical disk 624 such as a CD-ROM, DVD-ROM, or other  
2 optical media. The hard disk drive 616, magnetic disk drive 618, and optical disk  
3 drive 622 are each connected to the system bus 608 by one or more data media  
4 interfaces 626. Alternatively, the hard disk drive 616, magnetic disk drive 618,  
5 and optical disk drive 622 can be connected to the system bus 608 by one or more  
6 interfaces (not shown).

7 The disk drives and their associated computer-readable media provide non-  
8 volatile storage of computer readable instructions, data structures, program  
9 modules, and other data for computer 602. Although the example illustrates a hard  
10 disk 616, a removable magnetic disk 620, and a removable optical disk 624, it is to  
11 be appreciated that other types of computer readable media which can store data  
12 that is accessible by a computer, such as magnetic cassettes or other magnetic  
13 storage devices, flash memory cards, CD-ROM, digital versatile disks (DVD) or  
14 other optical storage, random access memories (RAM), read only memories  
15 (ROM), electrically erasable programmable read-only memory (EEPROM), and  
16 the like, can also be utilized to implement the exemplary computing system and  
17 environment.

18 Any number of program modules can be stored on the hard disk 616,  
19 magnetic disk 620, optical disk 624, ROM 612, and/or RAM 610, including by  
20 way of example, an operating system 626, one or more application programs 628,  
21 other program modules 630, and program data 632. Each of such operating  
22 system 626, one or more application programs 628, other program modules 630,  
23 and program data 632 (or some combination thereof) may implement all or part of  
24 the resident components that support the distributed file system.  
25

1 A user can enter commands and information into computer 602 via input  
2 devices such as a keyboard 634 and a pointing device 636 (e.g., a "mouse").  
3 Other input devices 638 (not shown specifically) may include a microphone,  
4 joystick, game pad, satellite dish, serial port, scanner, and/or the like. These and  
5 other input devices are connected to the processing unit 604 via input/output  
6 interfaces 640 that are coupled to the system bus 608, but may be connected by  
7 other interface and bus structures, such as a parallel port, game port, or a universal  
8 serial bus (USB).

9 A monitor 642 or other type of display device can also be connected to the  
10 system bus 608 via an interface, such as a video adapter 644. In addition to the  
11 monitor 642, other output peripheral devices can include components such as  
12 speakers (not shown) and a printer 646 which can be connected to computer 602  
13 via the input/output interfaces 640.

14 Computer 602 can operate in a networked environment using logical  
15 connections to one or more remote computers, such as a remote computing device  
16 648. By way of example, the remote computing device 648 can be a personal  
17 computer, portable computer, a server, a router, a network computer, a peer device  
18 or other common network node, and the like. The remote computing device 648 is  
19 illustrated as a portable computer that can include many or all of the elements and  
20 features described herein relative to computer 602.

21 Logical connections between computer 602 and the remote computer 648  
22 are depicted as a local area network (LAN) 650 and a general wide area network  
23 (WAN) 652. Such networking environments are commonplace in offices,  
24 enterprise-wide computer networks, intranets, and the Internet.  
25

1 When implemented in a LAN networking environment, the computer 602 is  
2 connected to a local network 650 via a network interface or adapter 654. When  
3 implemented in a WAN networking environment, the computer 602 typically  
4 includes a modem 656 or other means for establishing communications over the  
5 wide network 652. The modem 656, which can be internal or external to computer  
6 602, can be connected to the system bus 608 via the input/output interfaces 640 or  
7 other appropriate mechanisms. It is to be appreciated that the illustrated network  
8 connections are exemplary and that other means of establishing communication  
9 link(s) between the computers 602 and 648 can be employed.

10 In a networked environment, such as that illustrated with computing  
11 environment 600, program modules depicted relative to the computer 602, or  
12 portions thereof, may be stored in a remote memory storage device. By way of  
13 example, remote application programs 658 reside on a memory device of remote  
14 computer 648. For purposes of illustration, application programs and other  
15 executable program components such as the operating system are illustrated herein  
16 as discrete blocks, although it is recognized that such programs and components  
17 reside at various times in different storage components of the computing device  
18 602, and are executed by the data processor(s) of the computer.

19 Computer 602 typically includes at least some form of computer readable  
20 media. Computer readable media can be any available media that can be accessed  
21 by computer 602. By way of example, and not limitation, computer readable  
22 media may comprise computer storage media and communication media.  
23 Computer storage media includes volatile and nonvolatile, removable and non-  
24 removable media implemented in any method or technology for storage of  
25 information such as computer readable instructions, data structures, program

1 modules or other data. Computer storage media includes, but is not limited to,  
2 RAM, ROM, EEPROM, flash memory or other memory technology, CD-ROM,  
3 digital versatile disks (DVD) or other optical storage, magnetic cassettes, magnetic  
4 tape, magnetic disk storage or other magnetic storage devices, or any other media  
5 which can be used to store the desired information and which can be accessed by  
6 computer 602. Communication media typically embodies computer readable  
7 instructions, data structures, program modules or other data in a modulated data  
8 signal such as a carrier wave or other transport mechanism and includes any  
9 information delivery media. The term "modulated data signal" means a signal that  
10 has one or more of its characteristics set or changed in such a manner as to encode  
11 information in the signal. By way of example, and not limitation, communication  
12 media includes wired media such as wired network or direct-wired connection,  
13 and wireless media such as acoustic, RF, infrared and other wireless media.  
14 Combinations of any of the above should also be included within the scope of  
15 computer readable media.

16 The invention has been described herein in part in the general context of  
17 computer-executable instructions, such as program modules, executed by one or  
18 more computers or other devices. Generally, program modules include routines,  
19 programs, objects, components, data structures, etc. that perform particular tasks  
20 or implement particular abstract data types. Typically the functionality of the  
21 program modules may be combined or distributed as desired in various  
22 embodiments.

23 For purposes of illustration, programs and other executable program  
24 components such as the operating system are illustrated herein as discrete blocks,  
25 although it is recognized that such programs and components reside at various



1 times in different storage components of the computer, and are executed by the  
2 data processor(s) of the computer.

3 Alternatively, the invention may be implemented in hardware or a  
4 combination of hardware, software, and/or firmware. For example, one or more  
5 application specific integrated circuits (ASICs) could be designed or programmed  
6 to carry out the invention.

7 It should be noted that, although discussed primarily herein with reference  
8 to a serverless distributed file system, the invention can be used in any file system  
9 in which it is desired to identify identical files across multiple computers. Thus,  
10 the invention can be used in other embodiments, such as, for example, those with  
11 one or more centralized file servers.

## 12 13 14 **Conclusion**

15 Although the description above uses language that is specific to structural  
16 features and/or methodological acts, it is to be understood that the invention  
17 defined in the appended claims is not limited to the specific features or acts  
18 described. Rather, the specific features and acts are disclosed as exemplary forms  
19 of implementing the invention.  
20  
21  
22  
23  
24  
25